

IDENTIFICATION OF RICE VARIETIES

White Paper

DSC, Bellevue University

Prashant Raghuwanshi

DSC680-T301 Applied Data Science (2225-1)

Professor Catie Williams

04/09/2022

Identification of Rice Varieties Using Artificial Intelligence Methods

Towards a real-time rice grains sorting system. Identification of vitreous durum rice kernels using ANN based on their morphological, color, wavelet, and gaborlet features.

Business Problem :

The modern “Ready to Eat” Food Processing industry is importing grains (rice) from various international grains distributors. Their produced packed foods mostly depend on quality and varieties of imported raw grains. For them, it is important to keep the grains quality check on high priority.

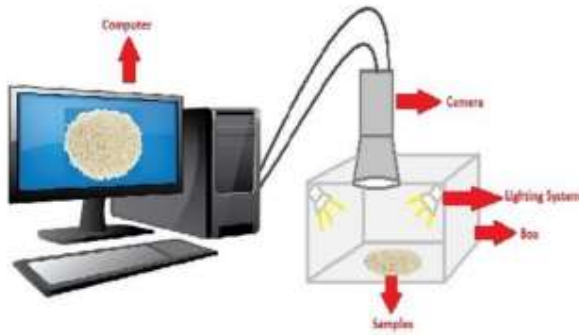
At present the quality check process completely depends on the manual steps performed by experienced grain experts and most of the time manual process is causing a delay in executing production lines.

Moreover, the manual quality check process result is not consistent and not quick, which required grain experts to be present on-site and near the production line.

The processing industry is looking for an Artificial intelligence-based grain sampling instrument that can quickly identify the grain quality and variety with maximum accuracy.

Background :

To solve the above-mentioned critical Business problem of “ready to eat” food processing industries. This project is going to develop the ML Model software which is going to be installed on a grain scanner machine and it will help the machine to identify the variety and quality of passed rice grain samples.



The images of rice samples were obtained first and images were processed using various image processing techniques. The resulting images are first converted to a grayscale image, then converted to a binary image and removed from the noise on the image. In the next phase, various morphological feature inference processes were applied to the obtained images. The classification phase of rice is given in Figure 1. During the modeling phase, the rice classification process was carried out using seven pieces of machine learning techniques. In the last step, the performance of the models used was evaluated.

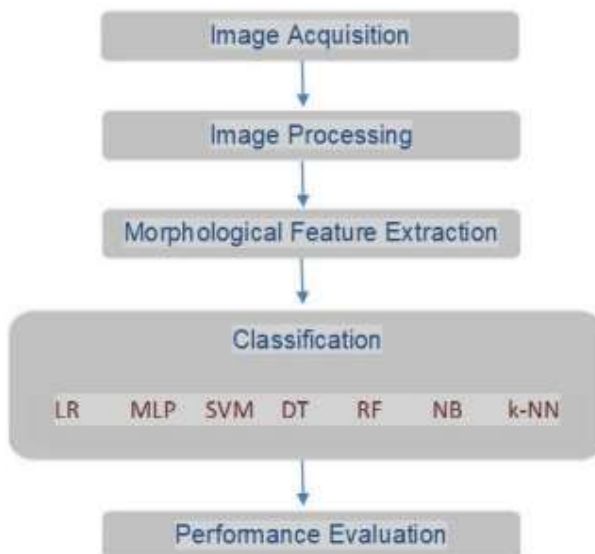
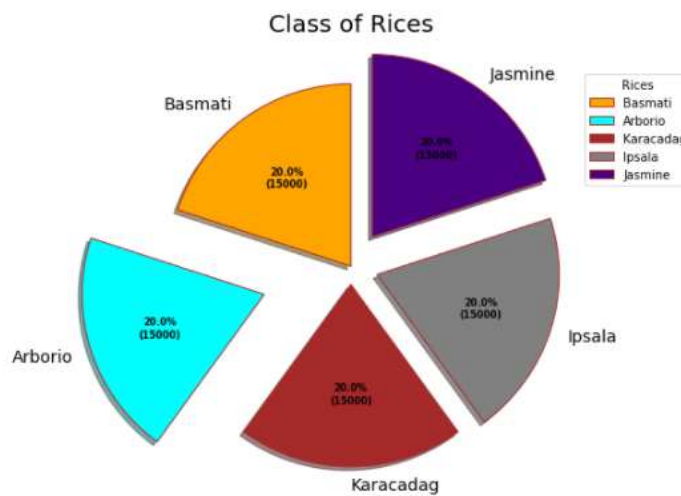


Fig. 1. Rice classification phases

Data Explanation:

A total of 75 thousand pieces of rice grain were obtained, including 15 thousand pieces of each variety of rice (Arborio, Basmati, Uppsala, Jasmine, Karacadag). Preprocessing operations were applied to the images and made available for feature extraction. A total of 106 features were inferred from the images; 12 morphological features and 4 shape features were obtained using morphological features and 90 color features were obtained from five different color spaces (RGB, HSV, Lab*, YCbCr, XYZ).



Attribute Information:

- 1.) Area: Returns the number of pixels within the boundaries of the rice grain.
- 2.) Perimeter: Calculates the circumference by calculating the distance between pixels around the boundaries of the rice grain.
- 3.) Major Axis Length: The longest line that can be drawn on the rice grain, i.e. the main axis distance, gives.
- 4.) Minor Axis Length: The shortest line that can be drawn on the rice grain, i.e. the small axis distance, gives.

5.) Eccentricity: It measures how round the ellipse, which has the same moments as the rice grain, is.

6.) Convex Area: Returns the pixel count of the smallest convex shell of the region formed by the rice grain.

7.) Extent: Returns the ratio of the region formed by the rice grain to the bounding box pixels.

8.) Class: Cammeo and Osmancik rice

Methods:

In this project, models were created using LR (Logistic Regression), MLP (Multilayer Perceptron), ABC (AdaBoostClassifier), DT (Decision Tree), RF (Random Forest), GB (Gradient Boosting), and kNN (K Nearest Neighbor) systems to classify rice grains according to their characteristics. (please refer to the code snippet in the appendix section)

Analysis :

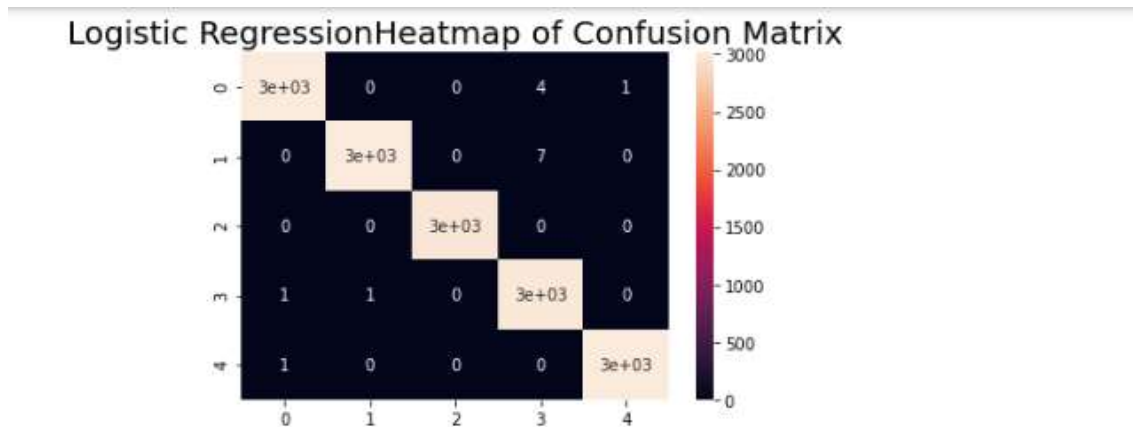
Classification performance measurement results are pasted below. all models except AdaBoostClassifier have achieved a classification success of over 99%. The 99.93% accuracy achieved in the LR model has the highest value among other models. Also for F1, recall and precision look good for the LR model.

Model Accuracy and F1 scores:

Logistic Regression Accuracy:	99.90%	F1-Score:	0.99900
Decision Tree Accuracy:	99.57%	F1-Score:	0.99567
Neural Network Accuracy:	99.93%	F1-Score:	0.99927
Random Forest Accuracy:	99.86%	F1-Score:	0.99860
Gradient Boosting Accuracy:	99.83%	F1-Score:	0.99833
AdaBoostClassifier Accuracy:	60.71%	F1-Score:	0.48542
KNeighborsClassifier Accuracy:	99.83%	F1-Score:	0.99827

Logistic Regressionclassification_report

	precision	recall	f1-score	support
Arborio	1.00	1.00	1.00	3010
Basmati	1.00	1.00	1.00	3022
Ipsala	1.00	1.00	1.00	2975
Jasmine	1.00	1.00	1.00	2985
Karacadag	1.00	1.00	1.00	3008
accuracy			1.00	15000
macro avg	1.00	1.00	1.00	15000
weighted avg	1.00	1.00	1.00	15000



Conclusion:

The logistic regression model is the best fit for rice gains predictions. We can use this classification model to develop Automatic systems that can be designed for many processes such as calibration of rice types and the separation of species from unwanted substances that may be present.

Assumption: Here I am assuming the Preprocessing data extraction operations is not having any limitations and the rice grains images were extracted & preprocessed successfully by image to data conversion tool.

Limitations: This Model is limited to identifying the Rice grains only. This model is suitable to identify Arborio, Basmati, Uppsala, Jasmine, and Karacadag rice varieties and would consider other varieties as unknown gain.

Challenges:

Major Constraints are related to used datasets and processed images, here the used datasets contain a total of 75 thousand rice grain images, including 15 thousand for each variety however due to the rapidly advancing Seed development process, we might not have all full collections of grain records under each gain varieties

Future Uses:

- 1) By adding additional grains datasets, we can increase the scope of this model to identify large varieties of grains.
- 2) Programmed in robot scanner machines and enable robots to filter adulterated gains

Recommendations:

Automatic systems can be designed for many processes such as calibration of rice types and the separation of species from unwanted substances that may be present. To increase the success rate in classification, more images can be obtained from species and it is thought that success rates can be increased by using morphological features as well as color and shape features.

Implementation Plan:

To develop ML applications, we have below implementation plan:

IDE's to use: Anaconda Jupyter

Coding languages: Python

Platform: AWS and local PC

Resource Data Science Engineer

Timeline: 03/15/2022 to 4/16/2022

Ethical Considerations:

- This Data contains processed physical images information related to multiple varieties of rice and does not contain any PII-related information.
- Datasets and information on data were extracted from the public websites → UCL machine learning repositories.
- This data research is not going to harm any privacy.

Appendix :

Data set Screenprint :

In [9]: `df.head(5)`

Out[9]:

	AREA	PERIMETER	MAJOR_AXIS	MINOR_AXIS	ECCENTRICITY	EQDIASQ	SOLIDITY	CONVEX_AREA	EXTENT	ASPECT_RATIO	ROUNDNESS	COMPACTNESS
0	7805	437.915	209.8215	48.0221	0.9735	99.6877	0.9775	7985	0.3547	4.3693	0.5114	0.4
1	7503	340.757	138.3361	69.8417	0.8632	97.7400	0.9660	7767	0.6637	1.9807	0.8120	0.7
2	5124	314.617	141.9803	46.5784	0.9447	80.7718	0.9721	5271	0.4760	3.0482	0.6505	0.5
3	7990	437.085	201.4386	51.2245	0.9671	100.8622	0.9659	8272	0.6274	3.9325	0.5256	0.5
4	7433	342.893	140.3350	68.3927	0.8732	97.2830	0.9831	7561	0.6006	2.0519	0.7944	0.6

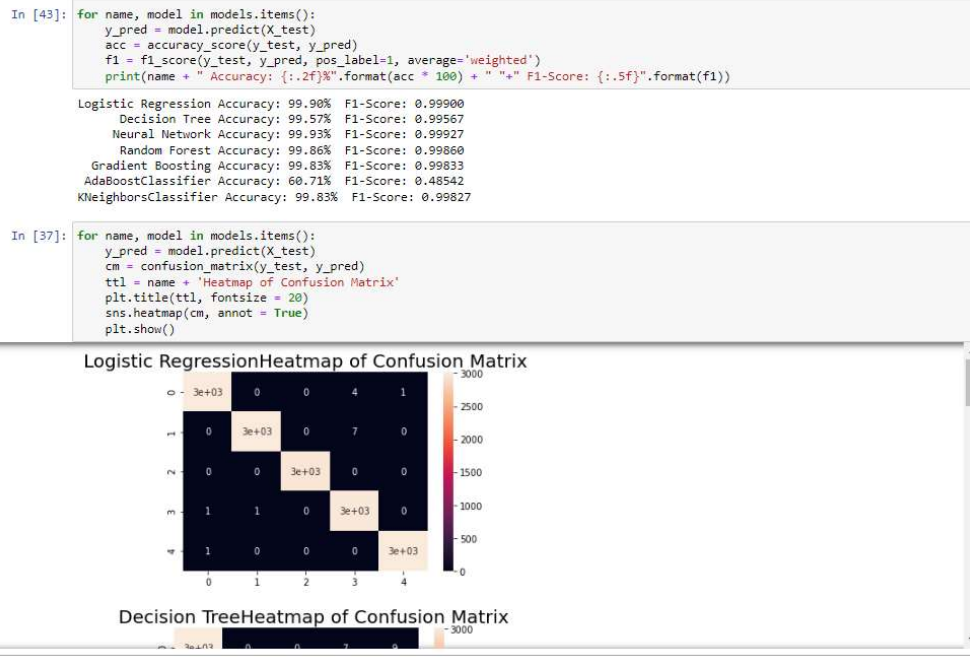
Model build screenprint:

```
In [32]: models = {
    "Logistic Regression": LogisticRegression(),
    "Decision Tree": DecisionTreeClassifier(),
    "Neural Network": MLPClassifier(),
    "Random Forest": RandomForestClassifier(),
    "Gradient Boosting": GradientBoostingClassifier(),
    "AdaBoostClassifier": AdaBoostClassifier(),
    "KNeighborsClassifier": KNeighborsClassifier()
}

for name, model in models.items():
    model.fit(X_train, y_train)
    print(name + " trained.")

Logistic Regression trained.
Decision Tree trained.
Neural Network trained.
Random Forest trained.
Gradient Boosting trained.
AdaBoostClassifier trained.
KNeighborsClassifier trained.
```

Accuracy computation :



Questions May be asked by the audience:

- 1) In the Real world how the data science team is going to capture the required rice grain data.

Answer: We have two ways to get rice grain preprocessed data

- Directly procure the required preprocessed grain data from agriculture university (or any available third party vendors) and have a contract with them for futures updates in rice varieties data
- Do the data processing in-house and procure the required equipment & Software for capturing & process the grain image data.
 - The images of rice samples were obtained first and images were processed using various image processing techniques. The resulting images are first converted to a grayscale image, then converted to a

binary image and removed from the noise on the image. In the next phase, various morphological feature inference processes were applied to the obtained images.

- 2) For procuring data from vendors or universities what ethical considerations do we need to take care of.

Answers:

- we are going to request data validity, original work certification, and copyright approval from vendors. Also need to make sure the vendor should be responsible to arrange the required government clearances.

- 3) What would be our revenue model?

Answer:

- we are going to provide installable software packages and their license to robotics scanner companies. And also having a plan to get some revenue from after-sales support from scanner companies.

- 4) The presented solution or model is meant for identifying five varieties of rice
How you are going to make it fit for other varieties?

Answer:

- our developed solution can be used for any additional add-on varieties of rice. However, before making it ready for any new variety of rice we need to retrain our model with updated or new datasets.

- 5) How do you conclude the LR model is the best fit?

Answer :

- We have built the model by using multiple classification algorithms and after computing the respective models' accuracy and other factors, we concluded the LR is going to best fit for our use case.

6) How accurate is your Model?

Answer:

- In this project, I have used clean and preprocessed data that contains equal no of samples counts for each variety. I believe for this project I have used only 5 varieties of rice grains and have got the model max accuracy of 99.8%.

7) With the increase in counts of varieties of grains, are you expecting any degradation in, accuracy?

Answer :

- Theoretically, we might get degradation inaccuracy, however, we are going to build and evaluate the models by using evolving multiple algorithms, which may result in a steady accuracy percentage (similar to current accuracy)

8) How the Model is designed to behave with unknown grain data.

Answer:

- At present, the model is not enough mature to identify the unknown gain variety data, but definitely, we are going to add the additional unknown classification variety in our training datasets.

9) please elaborate more on provided Recommendations:

Answer :

- This model is suitable for implementation in any Automatic gain scanning system and it can be designed for many gains identification processes such as calibration of rice types and the separation of species from unwanted substances that may be present.

10) Do you have a recommendation for making the model prediction more consistent:

Answer:

- To increase the success rate in classification, more images can be obtained from species and it is thought that success rates can be increased by using morphological features as well as color and shape features.

References :

<https://archive.ics.uci.edu/ml/datasets/Rice+%28Cammeo+and+Osmancik%29>

<https://www.muratkoklu.com/datasets/>

1: KOKLU, M., CINAR, I. and TASPINAR, Y. S. (2021). Classification of rice varieties with deep learning methods. Computers and Electronics in Agriculture, 187, 106285.

DOI: <https://doi.org/10.1016/j.compag.2021.106285>

2: CINAR, I. and KOKLU, M. (2021). Determination of Effective and Specific Physical Features of Rice Varieties by Computer Vision In Exterior Quality Inspection. Selcuk Journal of Agriculture and Food Sciences, 35(3), 229-243.

DOI: <https://doi.org/10.15316/SJAFS.2021.252>

3: CINAR, I. and KOKLU, M. (2022). Identification of Rice Varieties Using Machine Learning Algorithms. Journal of Agricultural Sciences.

DOI: <https://doi.org/10.15832/ankutbd.862482>

4: CINAR, I. and KOKLU, M. (2019). Classification of Rice Varieties Using Artificial Intelligence Methods. International Journal of Intelligent Systems and Applications in Engineering, 7(3), 188-194.

DOI: <https://doi.org/10.18201/ijisae.2019355381>